

- 19 -

WHAT IS CLAIMED IS:

1. A cDNA microarray data correction system for  
correcting global and local distortions of microarray data  
5 more precisely and correcting measurement errors caused by  
a difference in sensitivity between fluorescent dyes,  
comprising:

an input device for inputting previously-adjusted  
gene expression intensity data, considering flag  
10 information indicating a removal of background noise and  
reliability of each spot;

a data standardization means for standardizing the  
gene expression intensity data by using grid-by-grid order  
statistics for the input gene expression intensity data and  
15 for transmitting the standardized gene expression intensity  
data;

first correction means for estimating a distortion  
depending on a spot position on grid coordinates for the  
standardized gene expression intensity data by a  
20 nonparametric smoothing method and for transmitting first  
corrected gene expression intensity data whose distortion  
has been corrected; and

second correction means for performing an S-D  
transformation for the first corrected gene expression  
25 intensity data, for estimating a potential distortion  
caused by a difference in sensitivity between the  
fluorescent dyes in the gene expression intensity data by  
the nonparametric smoothing method, and for transmitting  
second corrected gene expression intensity data whose

- 20 -

distortion caused by the difference in sensitivity between the fluorescent dyes has been corrected; and an output device for outputting the second corrected gene expression intensity data.

5

2. The cDNA microarray data correction system according to claim 1, further comprising S-D transformation means for quantifying the distortion of the gene expression intensity data in an arbitrary stage and for visualizing it 10 on an S-D plot.

3. The cDNA microarray data correction system according to claim 1 or 2, wherein the order statistics are represented by the following EQ12 (where  $w_{ij}^k(c)$  is the 15 standardized gene expression intensity data,  $y_{ij}^k(c)$  is gene expression intensity data of all spots obtained in a channel, and  $L_k(c)$  and  $M_k(c)$  indicate 25 and 50 percent points of the gene expression intensity data obtained in channel  $c$  in grid  $k$ , respectively):

$$w_{ij}^k(c) = \frac{y_{ij}^k(c) - L_k(c)}{M_k(c) - L_k(c)}, \quad c=1,2, i=1,\dots,I, j=1,\dots,J, k=1,\dots,K. \quad (12)$$

4. The cDNA microarray data correction system according to claim 1 or 2, wherein the order statistics are represented by the following EQ13 (where  $w_{ij}^k(c)$  is the 25 standardized gene expression intensity data,  $y_{ij}^k(c)$  is gene expression intensity data of all spots obtained in a channel, and  $A_k(c)$ ,  $L_k(c)$  and  $M_k(c)$  indicate 35, 10 and 90 percent points of the gene expression intensity data

- 21 -

obtained in channel c in grid k, respectively):

$$w_{ij}^k(c) = \frac{y_{ij}^k(c) - A_k(c)}{M_k(c) - L_k(c)},$$

c = 1,2,i = 1,⋯,I,j = 1,⋯,J,k = 1,⋯,K. (13)

5        5. The cDNA microarray data correction system according to claim 3 or 4, wherein said data standardization means determines whether the gene expression intensity data of all spots obtained in at least two gene expression intensity data channels has been  
10 standardized and continues it until the gene expression intensity data of all spots has been standardized.

15        6. The cDNA microarray data correction system according to claim 1, wherein the standardized gene expression intensity data is represented by a sum of a true gene intensity and a distortion depending on the spot position.

20        7. The cDNA microarray data correction system according to claim 1, wherein said first correction means describes the distortion depending on the spot position by means of a nonparametric regression model represented by a regression relation of distortions with an x-axis, a y-axis, and an interaction of the x- and y-axes ( $\alpha_k^{(c)}(i), \beta_k^{(c)}(j)$ , and  
25  $\gamma_k^{(c)}((i - m_i)(j - m_j))$ , respectively) and estimates the distortion depending on the spot position ( $\xi_{ij}^k(c)$ ) by the nonparametric smoothing method represented by the following EQ14:

- 22 -

$$\hat{\xi}_{ij}^k(c) = \hat{\alpha}_x^{(c)}(i) + \hat{\beta}_x^{(c)}(j) + \hat{\gamma}_x^{(c)}((i-m_i)(j-m_j)), c=1,2, i=1,\dots,I, j=1,\dots,J. \quad (14)$$

8. The cDNA microarray data correction system according to claim 7, wherein the distortion depending on  
5 the spot position is corrected according to the following EQ15 (where  $\hat{z}_{ij}^k(c)$  is corrected true gene expression intensity data):

$$\hat{z}_{ij}^k(c) = w_{ij}^k(c) - \hat{\xi}_{ij}^k(c) \quad (15)$$

10 9. The cDNA microarray data correction system according to claim 8, wherein the S-D transformation in said second correction means is performed according to the following EQ16:

$$\begin{aligned} u_{ij}^k &= \hat{z}_{ij}^k(1) + \hat{z}_{ij}^k(2) \\ v_{ij}^k &= \hat{z}_{ij}^k(1) - \hat{z}_{ij}^k(2) \end{aligned} \quad (16)$$

15 10. The cDNA microarray data correction system according to claim 9, wherein said second correction means describes the distortion by means of a nonparametric regression model represented by the following EQ17, estimates a measurement error caused by the difference in sensitivity between the fluorescent dyes by a nonparametric smoothing method represented by the following EQ18 and EQ19, and corrects the error:

$$v_{ij}^k = \Phi(u_{ij}^k) + \epsilon_{ij}^k, \epsilon_{ij}^k \sim N(0, v^2) \quad (17)$$

- 23 -

$$\eta_{ij}^k = v_{ij}^k - \hat{\phi}(u_{ij}^k) \quad (18)$$

$$\begin{aligned}\hat{y}_{ij}^k(1) &= \frac{1}{2} (u_{ij}^k + \eta_{ij}^k) \\ \hat{y}_{ij}^k(2) &= \frac{1}{2} (u_{ij}^k - \eta_{ij}^k)\end{aligned} \quad (19)$$

11. The cDNA microarray data correction system

5 according to claim 1, wherein, supposing that a probability of gene expression is lower than 0.5, it is assumed for the correction that the fluorescence intensity detected at more than half of the spots within each grid indicates a background noise or a systematic error.

10

12. The cDNA microarray data correction system according to claim 11, wherein, supposing that  $L_k(c)$  and  $M_k(c)$  indicate 25 and 50 percent points of the fluorescence intensity obtained in at least two gene expression 15 intensity data channels in a grid, it is further assumed for the correction that  $L_k(c)$  and  $M_k(c) - L_k(c)$  are equal among the grids and teh channels on condition that most genes are in a non-expression state and that a distribution of 50 percent point or lower of the fluorescence intensity 20 is common to all grids and channels.

13. A cDNA microarray data correction method of correcting global and local distortions of microarray data more precisely and correcting measurement errors caused by 25 a difference in sensitivity between fluorescent dyes,

- 24 -

comprising the steps of:

inputting previously-adjusted gene expression intensity data, considering flag information indicating a removal of background noise and reliability of each spot;

5 standardizing the gene expression intensity data by using grid-by-grid order statistics for the input gene expression intensity data on condition that most genes are in a non-expression state;

outputting the standardized gene expression  
10 intensity data;

estimating a distortion depending on the spot position on grid coordinates for the standardized gene expression intensity data by a nonparametric smoothing method and correcting the data distortion depending on the  
15 spot position;

outputting the first corrected gene expression intensity data whose distortion depending on the spot position has been corrected;

performing an S-D transformation for the first  
20 corrected gene expression intensity data, estimating a potential distortion caused by a difference in sensitivity between the fluorescent dyes in the gene expression intensity data by the nonparametric smoothing method, and correcting the distortion caused by the difference in  
25 sensitivity between the fluorescent dyes; and

outputting the second corrected gene expression intensity data whose distortion caused by the difference in sensitivity between the fluorescent dyes has been corrected.

- 25 -

14. The cDNA microarray data correction method according to claim 13, further comprising a step of quantifying the distortion of the gene expression intensity data in an arbitrary stage and visualizing it on an S-D plot.

15. The cDNA microarray data correction method according to claim 13 or 14, wherein the order statistics are represented by the following EQ20 (where  $w_{ij}^k(c)$  is the 10 standardized gene expression intensity data,  $y_{ij}^k(c)$  is gene expression intensity data of all spots obtained in a channel, and  $L_k(c)$  and  $M_k(c)$  indicate 25 and 50 percent points of the gene expression intensity data obtained in channel c in grid k, respectively):

$$w_{ij}^k(c) = \frac{y_{ij}^k(c) - L_k(c)}{M_k(c) - L_k(c)}, \quad c=1,2, i=1, \dots, I, j=1, \dots, J, k=1, \dots, K. \quad (20)$$

16. The cDNA microarray data correction method according to claim 13 or 14, wherein the order statistics are represented by the following EQ21 (where  $w_{ij}^k(c)$  is the 20 standardized gene expression intensity data,  $y_{ij}^k(c)$  is gene expression intensity data of all spots obtained in a channel, and  $A_k(c)$ ,  $L_k(c)$  and  $M_k(c)$  indicate 35, 10 and 90 percent points of the gene expression intensity data obtained in channel c in grid k, respectively):

$$w_{ij}^k(c) = \frac{y_{ij}^k(c) - A_k(c)}{M_k(c) - L_k(c)}, \\ c = 1,2, i = 1, \dots, I, j = 1, \dots, J, k = 1, \dots, K. \quad (21)$$

- 26 -

17. The cDNA microarray data correction method according to claim 15 or 16, wherein, in the step of standardizing the data, it is determined whether the gene expression intensity data of all spots obtained in at least 5 two gene expression intensity data channels have been standardized and it is continued until the gene expression intensity data of all spots have been standardized.

18. The cDNA microarray data correction method according to claim 17, wherein the standardized gene 10 expression intensity data is represented by a sum of a true gene intensity and a distortion depending on the spot position.

19. The cDNA microarray data correction method according to claim 13, wherein, in the step of correcting 15 the data distortion depending on the spot position, the distortion depending on the spot position is described by means of a nonparametric regression model represented by a regression relation of distortions with an x-axis, a y-axis, and an interaction of the x- and y-axes ( $\alpha_k^{(c)}(i)$ ,  $\beta_k^{(c)}(j)$ , and 20  $\gamma_k^{(c)}((i - m_i)(j - m_j))$ , respectively) and the distortion depending on the spot position ( $\xi_{ij}^k(c)$ ) is estimated by the nonparametric smoothing method represented by the following EQ22:

$$\begin{aligned} \hat{\xi}_{ij}^k(c) &= \hat{\alpha}_k^{(c)}(i) + \hat{\beta}_k^{(c)}(j) + \hat{\gamma}_k^{(c)}((i - m_i)(j - m_j)), \\ c &= 1, 2, i = 1, \dots, I, j = 1, \dots, J. \quad (22) \end{aligned}$$

20. The cDNA microarray data correction method according to claim 19, wherein the distortion depending on

- 27 -

the spot position is corrected according to the following EQ23 (wherein  $\hat{z}_{ij}^k(c)$  is corrected true gene expression intensity data):

$$\hat{z}_{ij}^k(c) = w_{ij}^k(c) - \xi_{ij}^k(c) \quad (23)$$

5

21. The cDNA microarray data correction method according to claim 19, wherein the S-D transformation in the step of correcting the distortion caused by the difference in sensitivity between fluorescent dyes is performed according to the following EQ24:

$$\begin{aligned} u_{ij}^k &= \hat{z}_{ij}^k(1) + \hat{z}_{ij}^k(2) \\ v_{ij}^k &= \hat{z}_{ij}^k(1) - \hat{z}_{ij}^k(2) \end{aligned} \quad (24)$$

22. The cDNA microarray data correction method according to claim 20, wherein, in the step of correcting the distortion caused by the difference in sensitivity between the fluorescent dyes, the distortion is described by means of a nonparametric regression model represented by the following EQ25, a measurement error caused by the difference in sensitivity between the fluorescent dyes is estimated by a nonparametric smoothing method represented by the following EQ26 and EQ27, and the error is corrected:

$$v_{ij}^k = \phi(u_{ij}^k) + \varepsilon_{ij}^k, \quad \varepsilon_{ij}^k = N(0, v^2) \quad (25)$$

$$\eta_{ij}^k = v_{ij}^k - \hat{\phi}(u_{ij}^k) \quad (26)$$

- 28 -

$$\begin{aligned}\hat{y}_{ij}^k(1) &= \frac{1}{2} (u_{ij}^k + \eta_{ij}^k) \\ \hat{y}_{ij}^k(2) &= \frac{1}{2} (u_{ij}^k - \eta_{ij}^k)\end{aligned}\quad (27)$$

23. The cDNA microarray data correction method according to claim 13, wherein, supposing that a  
5 probability of gene expression is lower than 0.5, it is assumed for the correction that the fluorescence intensity detected at more than half of the spots within each grid indicates a background noise or a systematic error.

10 24. The cDNA microarray data correction method according to claim 23, wherein, supposing that  $L_k(c)$  and  $M_k(c)$  indicate 25 and 50 percent points of the fluorescence intensity obtained in at least two gene expression intensity data channels in a grid, it is further assumed  
15 for the correction that  $L_k(c)$  and  $M_k(c) - L_k(c)$  are equal among the grids and the channels on condition that most genes are in a non-expression state and that a distribution of 50 percent point or lower of the fluorescence intensity is common to all grids and channels.

20 25. The cDNA microarray data correction method according to claim 23, wherein, denoting that  $A_k(c)$ ,  $L_k(c)$  and  $M_k(c)$  indicate 35, 10 and 90 percent points of the fluorescence in a grid  $k$  for channel  $c$ , it is assumed for  
25 the correction that  $A_k(c)$  and  $M_k(c) - L_k(c)$  are common to all grids and channels.

- 29 -

26. A cDNA microarray data correction program for use in correcting global and local distortions of microarray data more precisely and correcting measurement errors caused by a difference in sensitivity between

5 fluorescent dyes with a computer to execute the steps of:

inputting previously-adjusted gene expression intensity data, considering flag information indicating a removal of background noise and reliability of each spot;

standardizing the gene expression intensity data  
10 by using grid-by-grid order statistics for the input gene expression intensity data on condition that most genes are in a non-expression state;

outputting the standardized gene expression intensity data;

15 estimating a distortion depending on the spot position on grid coordinates for the standardized gene expression intensity data by a nonparametric smoothing method and correcting the data distortion depending on the spot position;

20 outputting the first corrected gene expression intensity data whose distortion depending on the spot position has been corrected;

25 performing an S-D transformation for the first corrected gene expression intensity data, estimating a potential distortion caused by a difference in sensitivity between the fluorescent dyes in the gene expression intensity data by the nonparametric smoothing method, and correcting the distortion caused by the difference in sensitivity between the fluorescent dyes; and

- 30 -

outputting the second corrected gene expression intensity data whose distortion caused by the difference in sensitivity between the fluorescent dyes has been corrected.

- 5        27. A computer-readable memory medium containing a cDNA microarray data correction program for use in correcting global and local distortions of microarray data more precisely and correcting measurement errors caused by a difference in sensitivity between  
10      fluorescent dyes with a computer to execute the steps of:  
              inputting previously-adjusted gene expression intensity data, considering flag information indicating a removal of background noise and reliability of each spot;  
              standardizing the gene expression intensity data  
15      by using grid-by-grid order statistics for the input gene expression intensity data on condition that most genes are in a non-expression state;  
              outputting the standardized gene expression intensity data;  
20      estimating a distortion depending on the spot position on grid coordinates for the standardized gene expression intensity data by a nonparametric smoothing method and correcting the data distortion depending on the spot position;  
25      outputting the first corrected gene expression intensity data whose distortion depending on the spot position has been corrected;  
              performing an S-D transformation for the first corrected gene expression intensity data, estimating a

- 31 -

potential distortion caused by a difference in sensitivity  
between the fluorescent dyes in the gene expression  
intensity data by the nonparametric smoothing method, and  
correcting the distortion caused by the difference in  
5 sensitivity between the fluorescent dyes; and  
outputting the second corrected gene expression  
intensity data whose distortion caused by the difference in  
sensitivity between the fluorescent dyes has been corrected.